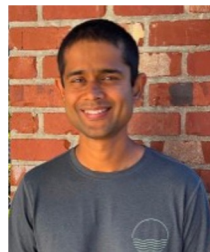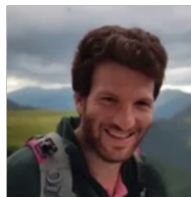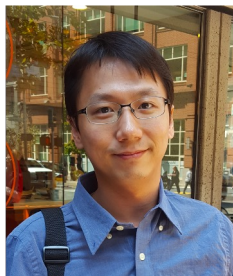# ManagerTower: Aggregating the Insights of Uni-Modal Experts for Vision-Language Representation Learning

Xiao Xu[1,3], Bei Li[2,3], Chenfei Wu[3], Shao-Yen Tseng[4], Anahita Bhiwandiwalla[4], Shachar Rosenman[4], Vasudev Lal[4], Wanxiang Che[1], Nan Duan[3]

[1]Harbin Institute of Technology, [2]Northeastern University, [3]Microsoft Research Asia, [4]Intel Labs

ACL 2023 Oral   ||   Presenter: Xiao Xu   ||   July, 2023

1

# Outline

- Background
- Motivation
- Architecture & Manager Design
- Visualization

# Background: Vision-Language Learning
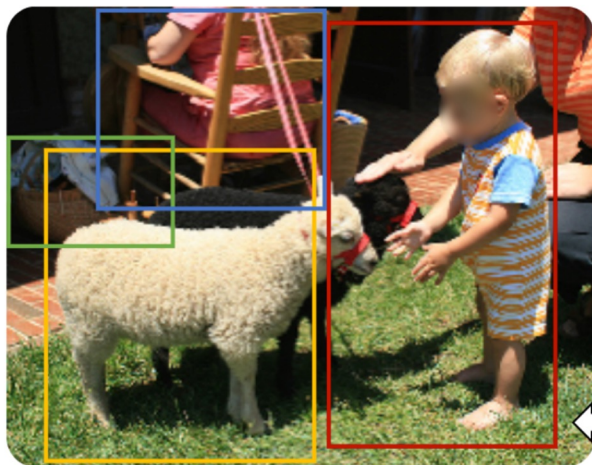
# What is Vision-Language (VL) Learning?

Goal: Train a smart AI system that can understand both image and text.

Approach: Large-scale self-supervised pre-training on image-text pairs.

Image from: https://arxiv.org/abs/1912.02315.

# Two-Tower Architecture



Cross-Modal Encoder

Textual Encoder

Visual Encoder

Two-Tower

# Two-Tower vs. BridgeTower

Two-Tower

BridgeTower

# Motivation: Adaptively Exploit Uni-Modal Insights

# Limitations of BridgeTower

- Ineffective layer-by-layer utilization
- The number of cross-modal layers is tied to the number of uni-modal layer representations it used



(a) BridgeTower

# BridgeTower vs. ManagerTower

## Limitations of BridgeTower

- Ineffective layer-by-layer utilization
- The number of cross-modal layers is tied to the number of uni-modal layer representations it used

## Advances of ManagerTower

- Takes multi-layer uni-modal representations as the insights of pre-trained uni-modal experts at different levels
- Adaptively aggregates insights via managers in each cross-modal layer



Textual Encoder    Cross-Modal Encoder    Visual Encoder

(a) BridgeTower

(b) ManagerTower

# Architecture & Manager Design

# ManagerTower Architecture



$\mathbf{C}_\ell^T$

$\mathbf{C}_\ell^V$

$\times 6$

Feed Forward

Cross Attention

Self Attention

Feed Forward

Cross Attention

Self Attention

$\mathbf{T} = [\mathbf{T}_7, \dots, \mathbf{T}_{12}]$

$\mathbf{V} = [\mathbf{V}_7, \dots, \mathbf{V}_{12}]$

Feed Forward

$\times 12$

Self Attention

Feed Forward

$12 \times$

Self Attention

Manager

Manager

Managers and Experts.

$\mathbf{C}_{\ell-1}^T$

$\mathbf{C}_{\ell-1}^V$

Textual Encoder

Cross-Modal Encoder

Visual Encoder

ManagerTower can work with any visual, textual, or cross-modal encoder.

# Static Aggregation of Experts (SAE) Manager



$\ell$: cross-modal layer index

$$\text{Input} \begin{cases} \text{Uni-Modal Part:} & \text{L} \times \text{D} \\ \text{Cross-Modal Part:} & \text{L} \times \text{D} \end{cases} \xrightarrow{\text{Bridge}} \text{Output} : \text{L} \times \text{D} \qquad (1)$$

$$\text{Input} \begin{cases} \text{Uni-Modal Part:} & 6 \times \text{L} \times \text{D} \\ \text{Cross-Modal Part:} & (\ell - 1) \times \text{L} \times \text{D} \\ \text{Learned Weights:} & (6 + \ell - 1) \times \text{D} \end{cases} \xrightarrow{\text{Manager}} \begin{cases} \text{Uni-Modal Aggregated:} & \text{L} \times \text{D} \\ \text{Cross-Modal Aggregated:} & \text{L} \times \text{D} \end{cases} \xrightarrow{\text{Manager}} \text{Output} : \text{L} \times \text{D} \quad (2)$$

**BridgeTower**

**SAE**

- Textual Expert
- Cross-Modal Expert
- Visual Expert

Cosine similarity of aggregated representations between every two consecutive managers

Textual Managers

Visual Managers

SAE-Uni
SAE-Cross

1. Uni-modal similarity ≈ 1
2. Cross-modal similarity increases with depth and gets closer to 1

# Static Aggregation of Uni-Modal Experts (SAUE) Manager



$$\text{Input} \begin{cases} \text{Uni-Modal Part:} & 6 \times L \times D \\ \text{Cross-Modal Part:} & (\ell - 1) \times L \times D \\ \text{Learned Weights:}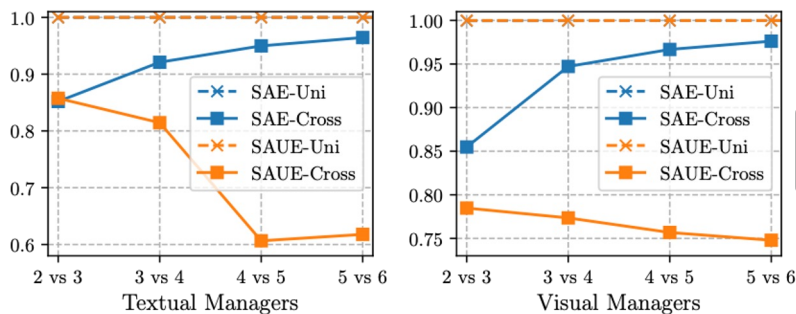 & (6 + \ell - 1) \times D \end{cases} \xrightarrow{\text{Manager}} \begin{cases} \text{Uni-Modal Aggregated:} & L \times D \\ \text{Cross-Modal Aggregated:} & L \times D \end{cases} \xrightarrow{\text{Manager}} \text{Output} : L \times D \quad (2)$$

$$\text{Input} \begin{cases} \text{Uni-Modal Part:} & 6 \times L \times D \\ \text{Cross-Modal Part:} & L \times D \\ \text{Learned Weights:} & 7 \times D \end{cases} \xrightarrow{\text{Manager}} \begin{cases} \text{Uni-Modal Aggregated:} & L \times D \\ \text{Cross-Modal Aggregated:} & L \times D \end{cases} \xrightarrow{\text{Manager}} \text{Output} : L \times D \quad (3)$$

Cross-modal similarity decreases with depth

1. Uni-modal similarity still ≈ 1
2. Input-independent learned weights: N × D

💡 Intuition: the need for uni-modal semantic knowledge varies among cross-modal layers, tokens and samples.

13

$$\text{Input} \begin{cases} \text{Uni-Modal Part:} & 6 \times L \times D \\ \text{Cross-Modal Part:} & L \times D \\ \boxed{\text{Learned Weights:} \quad 7 \times D} \end{cases} \xrightarrow{\text{Manager}} \begin{cases} \text{Uni-Modal Aggregated:} & L \times D \\ \text{Cross-Modal Aggregated:} & L \times D \end{cases} \xrightarrow{\text{Manager}} \text{Output}: L \times D \quad (3)$$

$$\text{Input} \begin{cases} \text{Uni-Modal Part:} & 6 \times L \times D \\ \text{Cross-Modal Part:} & L \times D \\ \boxed{\text{Generated Weights:} \quad 7 \times L} \end{cases} \xrightarrow{\text{Manager}} \begin{cases} \text{Uni-Modal Aggregated:} & L \times D \\ \text{Cross-Modal Aggregated:} & L \times D \end{cases} \xrightarrow{\text{Manager}} \text{Output}: L \times D \quad (4)$$
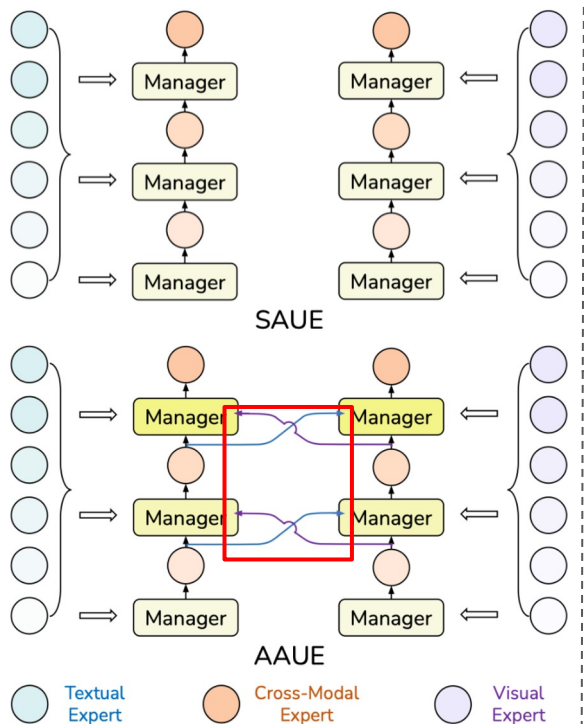
SAE & SAUE managers: learned weights, static sentence-level managers
AAUE managers: generated weights, adaptive token-level managers

| Type | Visual Query | Weight | Test-Dev | $R_{\text{MEAN}}$ |
|------|--------------|--------|----------|-------------------|
| BridgeTower | - | $N \times 1$ | 75.91 | 93.33 |
| SAE | - | $N \times 1$ | 76.19 | 93.57 |
|  | - | $N \times D$ | 76.18 | 93.73 |
| SAUE | - | $N \times 1$ | 76.38 | 93.75 |
|  | - | $N \times D$ | 76.55 | 93.82 |
| AAUE | $\mathbf{C}_{\ell-1}^{\text{V}}$ | $N \times L$ | 76.52 | 93.84 |
|  | $\boxed{\text{CA}(\mathbf{C}_{\ell-1}^{\text{V}}, \mathbf{C}_{\ell-1}^{\text{T}})}$ | $N \times L$ | **76.65** | **93.97** |

CA: Cross-Attention

Cross-Modal Fused Query: $\quad L_V \times D, L_T \times D \xrightarrow{\text{Cross-Attention}} L_V \times D$

AAUE managers achieves best performance.

# Main Results

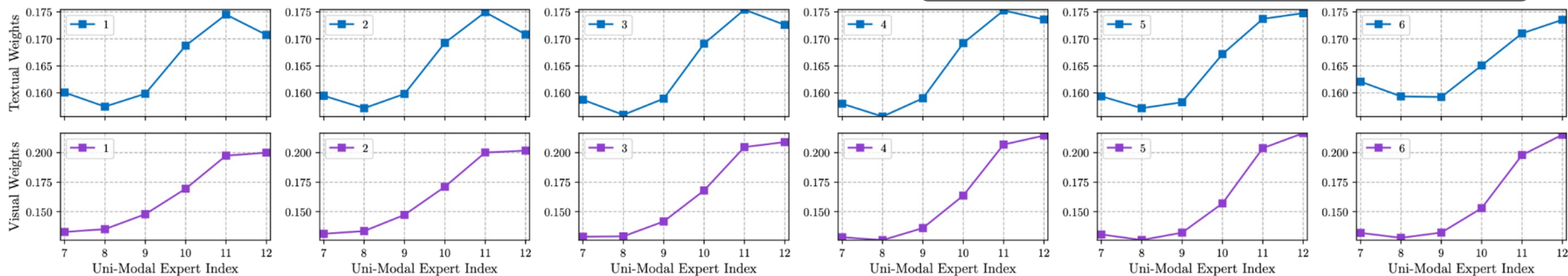| Model | # Pre-train Images | Visual Backbone | VQAv2 Test-Dev | VQAv2 Test-Std | SNLI-VE Dev | SNLI-VE Test | NLVR$^2$ Dev | NLVR$^2$ Test-P | Flickr30K IR@1 | Flickr30K TR@1 |
|---|---|---|---|---|---|---|---|---|---|---|
| *Base-size models pre-trained on 4M public data* | | | | | | | | | | |
| ViLT$_{BASE}$ (Kim et al., 2021) | 4M | ViT-B-384/32 | 71.26 | - | - | - | 75.70 | 76.13 | 64.4 | 83.5 |
| UNITER$_{BASE}$ (Chen et al., 2020) ∗ | 4M | Faster R-CNN | 72.70 | 72.91 | 78.59 | 78.28 | 77.18 | 77.85 | 72.52 | 85.90 |
| VILLA$_{BASE}$ (Gan et al., 2020) ∗ | 4M | Faster R-CNN | 73.59 | 73.67 | 79.47 | 79.03 | 78.39 | 79.30 | 74.74 | 86.60 |
| UNIMO$_{BASE}$ (Li et al., 2021b) | 4M | Faster R-CNN | 73.79 | 74.02 | 80.00 | 79.10 | - | - | 74.66 | 89.70 |
| ALBEF$_{BASE}$ (Li et al., 2021a) ∗ | 4M | DeiT-B-224/16 | 74.54 | 74.70 | 80.14 | 80.30 | 80.24 | 80.50 | 82.8 | 94.3 |
| VinVL$_{BASE}$ (Zhang et al., 2021) | 5.7M | ResNeXt-152 | 75.95 | 76.12 | - | - | 82.05 | 83.08 | - | - |
| METER-Swin$_{BASE}$ (Dou et al., 2022) | 4M | Swin-B-384/32 | 76.43 | 76.42 | 80.61 | 80.45 | 82.23 | 82.47 | 79.02 | 92.40 |
| VLMO$_{BASE}$ (Wang et al., 2021a) | 4M | BEiT-B-224/16 | 76.64 | 76.89 | - | - | 82.77 | 83.34 | 79.3 | 92.3 |
| METER-CLIP$_{BASE}$ (Dou et al., 2022) | 4M | CLIP-ViT-B-224/16 | 77.68 | 77.64 | 80.86 | 81.19 | 82.33 | 83.05 | 82.22 | 94.30 |
| BridgeTower$_{BASE}$ (Xu et al., 2022) | 4M | CLIP-ViT-B-224/16 | 78.66 | 78.73 | 81.11 | 81.19 | 81.85 | 83.09 | 85.83 | 94.73 |
| ManagerTower$_{BASE}$ (**Ours**) | 4M | CLIP-ViT-B-224/16 | **79.39** | **79.15** | **81.26** | **81.44** | **82.81** | **83.34** | **86.56** | **95.64** |
| *Models pre-trained on more data and/or with larger size* | | | | | | | | | | |
| UNITER$_{LARGE}$ (Chen et al., 2020) ∗ | 4M | Faster R-CNN | 73.82 | 74.02 | 79.39 | 79.38 | 79.12 | 79.98 | 75.56 | 87.30 |
| VILLA$_{LARGE}$ (Gan et al., 2020) ∗ | 4M | Faster R-CNN | 74.69 | 74.87 | 80.18 | 80.02 | 79.76 | 81.47 | 76.26 | 87.90 |
| UNIMO$_{LARGE}$ (Li et al., 2021b) | 4M | Faster R-CNN | 75.06 | 75.27 | 81.11 | 80.63 | - | - | 78.04 | 89.40 |
| ALBEF$_{BASE}$ (Li et al., 2021a) ∗ | 14M | DeiT-B-224/16 | 75.84 | 76.04 | 80.80 | 80.91 | 82.55 | 83.14 | 85.6 | 95.9 |
| VinVL$_{LARGE}$ (Zhang et al., 2021) | 5.7M | ResNeXt-152 | 76.52 | 76.63 | - | - | 82.67 | 83.98 | - | - |
| BLIP$_{BASE}$ (Li et al., 2022a) ∗ | 14M | DeiT-B-224/16 | 77.54 | 77.62 | - | - | 82.67 | 82.30 | 87.2 | 96.6 |
| SimVLM$_{BASE}$ (Wang et al., 2021b) ⋆ | 1.8B | ResNet-101 | 77.87 | 78.14 | 84.20 | 84.15 | 81.72 | 81.77 | - | - |
| BLIP$_{BASE}$ (Li et al., 2022a) ∗ | 129M | DeiT-B-224/16 | 78.24 | 78.17 | - | - | 82.48 | 83.08 | 87.3 | 97.3 |
| SimVLM$_{LARGE}$ (Wang et al., 2021b) ⋆ | 1.8B | ResNet-152 | 79.32 | 79.56 | 85.68 | 85.62 | 84.13 | 84.84 | - | - |
| VLMO$_{LARGE}$ (Wang et al., 2021a) | 4M | BEiT-L-224/16 | 79.94 | 79.98 | - | - | 85.64 | 86.86 | 84.5 | 95.3 |
| SimVLM$_{HUGE}$ (Wang et al., 2021b) ⋆ | 1.8B | Larger ResNet-152 | 80.03 | 80.34 | 86.21 | 86.32 | 84.53 | 85.15 | - | - |

Follow METER's and BridgeTower's setting + 4M Vision-Language Pre-training + Managers => significant gains and outperforms some models trained with more data and parameters.

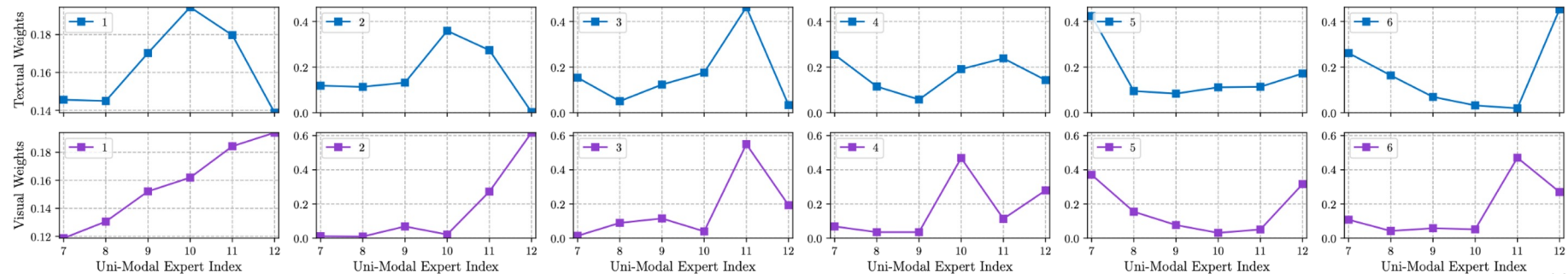# Visualization of Aggregation Weights



SAUE Managers — Horizontal: similar progressive weight distributions

AAUE Managers — Horizontal: diverse weight distributions

16

# Take-Away Messages

- Introduce managers in each cross-modal layer to
  - adaptively aggregate the insights of pre-trained uni-modal experts at different levels
  - flexibly generate different aggregation weights for different tokens in different samples
  - facilitate more comprehensive cross-modal alignment and fusion
- Cross-modal fused query
  - incorporates the output visual & textual representations of the previous cross-modal layer
  - to help managers to correctly aggregate uni-modal semantic knowledge required by the current cross-modal layer
- ManagerTower can work with any visual, textual, or cross-modal encoder

# Thanks & QA

Xiao Xu[1,3], Bei Li[2,3], Chenfei Wu[3], Shao-Yen Tseng[4], Anahita Bhiwandiwalla[4], Shachar Rosenman[4], Vasudev Lal[4], Wanxiang Che[1], Nan Duan[3]

[1]Harbin Institute of Technology, [2]Northeastern University, [3]Microsoft Research Asia, [4]Intel Labs

ACL 2023 Oral   ||   Presenter: Xiao Xu   ||   July, 2023

Slides and more in https://looperxx.github.io/.